



[biblio.ugent.be](http://biblio.ugent.be)

The UGent Institutional Repository is the electronic archiving and dissemination platform for all UGent research publications. Ghent University has implemented a mandate stipulating that all academic publications of UGent researchers should be deposited and archived in this repository. Except for items where current copyright restrictions apply, these papers are available in Open Access.

This item is the archived peer-reviewed author-version of: Robust calibrations on reduced sample sets for API content prediction in tablets: Definition of a cost-effective NIR model development strategy

Authors: Sigrid Pieters, Wouter Saeys, Tom Van den Kerkhof, Mohammad Goodarzi, Mario Hellings, Thomas De Beer, Yvan Vander Heyden

In: *Analytica Chimica Acta*, 761, 62-70 (2013)

Optional: link to the article

**To refer to or to cite this work, please use the citation to the published version:**

**Authors (year). Title. *journal* Volume(Issue) page-page. Doi** 10.1016/j.aca.2012.11.034

# **Robust calibrations on reduced sample sets for API content prediction in tablets: Definition of a cost-effective NIR model development strategy**

Sigrid Pieters<sup>1</sup>, Wouter Saeys<sup>2</sup>, Tom Van den Kerkhof<sup>3</sup>, Mohammad Goodarzi<sup>2</sup>, Mario Hellings<sup>3</sup>, Thomas De Beer<sup>4</sup>, Yvan Vander Heyden<sup>1</sup>

<sup>1</sup> Department of Analytical Chemistry and Pharmaceutical Technology, Center for Pharmaceutical Research, Vrije Universiteit Brussel - VUB, Laarbeeklaan 103, B-1090 Brussels, Belgium

<sup>2</sup> Department of Biosystems, Faculty of Bioscience Engineering, Katholieke Universiteit Leuven – KU Leuven, Kasteelpark Arenberg 30, B-3001 Heverlee, Belgium

<sup>3</sup> ChemPhar development Johnson & Johnson, Turnhoutseweg 30, B-2340 Beerse, Belgium

<sup>4</sup> Laboratory of Pharmaceutical Process Analytical Technology, Universiteit Gent, Harelbekestraat 72, B-9000 Ghent, Belgium

\* Corresponding author: Tel.: +32 2 477 47 34  
Fax: +32 2 477 47 35  
Email: [yvanvdh@vub.ac.be](mailto:yvanvdh@vub.ac.be) (Y. Vander Heyden)

Keywords: near infrared spectroscopy, prior spectral information, tablets, orthogonal projection methods, augmentation methods

## **Abstract**

Owing to spectral variations from other sources than the component of interest, large investments in the NIR model development may be required to obtain satisfactory and robust prediction performance. To make the NIR model development for routine active pharmaceutical ingredient (API) prediction in tablets more cost-effective, alternative modelling strategies were proposed. They used a massive amount of prior spectral information on intra- and inter-batch variation and the pure component spectra to define a clutter, i.e. the detrimental spectral information. This was subsequently used for artificial data augmentation and/or orthogonal projections. The model performance improved statistically significantly, with a 34-40% reduction in RMSEP while needing fewer model latent variables, by applying the following procedure before PLS regression: (1) augmentation of the calibration spectra with the spectral shapes from the clutter, and (2) net analyte pre-processing (NAP). The improved prediction performance was not compromised when reducing the variability in the calibration set, making exhaustive calibration unnecessary. Strong water content variations in the tablets caused frequency shifts of the API absorption signals that could not be

1 included in the clutter. Updating the model for this kind of variation demonstrated that the  
2 completeness of the clutter is critical for the performance of these models and that the model will  
3 only be more robust for spectral variation that is not co-linear with the one from the property of  
4 interest.

## 6 **1. Introduction**

8 Together with the growing popularity of near-infrared (NIR) spectroscopy, the interest to use  
9 multivariate prediction models (e.g. PLS) emerged in different industries (food, chemical,  
10 pharmaceutical...) [1-3]. Compared to HPLC, NIR spectroscopy provides a non-destructive, extremely  
11 fast, and hands-on analysis, which can substantially reduce analysis times and costs for the routine  
12 assay of active pharmaceutical ingredient (API) in tablets. Multivariate calibration models (e.g. PLS)  
13 are typically built in an empirical way. Although this may offer flexibility, it can also be a drawback  
14 making the model development, maintenance and update more complex [4]. Consequently, more  
15 investment may be required, which must guarantee that a low prediction error is obtained and  
16 maintained for future predictions of novel batches.

17 There are numerous examples where spectral disturbances from various sources, e.g.  
18 variations in a process, temperature, humidity... , upset the model predictions. In this paper, we  
19 specifically consider this problem in tablets produced from a routine production line. Besides  
20 measurement variability (path length, temperature...), systematic changes in the measured signals in  
21 the spectra of tablets may appear, e.g. originating from process and raw material variability (inter-  
22 and intra-batch effects). A tableting process involves different subsequent processing steps, such as  
23 mixing the ingredients, granulation, compaction and coating. Some tablet production variables, e.g.  
24 grain size, compaction pressure, and coating thickness, may introduce spectral variability that is  
25 irrelevant for predicting the API content in the tablets [5]. Also, the excipients in the tablet can vary  
26 slightly between different batches or samples (excipient homogeneity), while the tablet composition  
27 can vary as well when different API target concentration levels are considered in the model.  
28 Depending on the relative humidity conditions during their production, packaging, storage and  
29 analysis, the water content is another important variable that may upset the NIR model predictions  
30 [6].

31 It is evident that when more relevant variation is covered by the calibration set of the NIR  
32 model, lower prediction errors can be obtained on new independent samples [7]. Including all the  
33 expected variation in the calibration set is the most intuitive way for doing this, and is still  
34 recommended by regulatory authorities [8]. Drawbacks of this approach are that the model  
35 complexity may increase rapidly when more 'non-relevant' spectral variability is included in the

exhaustive calibration, needing additional latent variables or introducing non-linearities that may require non-linear calibration techniques. A more important practical problem is that a substantial number of well-chosen samples, describing all the expected variability, is needed to obtain robust model performance. It is a non-trivial task to obtain a good estimate of the expected variation in the samples to be predicted, also because there is a practical limit (particularly in terms of costs) to the number of samples to be analysed with the reference technique. Hence, in most practical cases one selects randomly the calibration samples from 'representative' batches. As a result, existing calibration bases rarely contain all the relevant variability from all the influence factors which can occur in industrial conditions [9], making the models less robust for their long term use and requiring frequently model updates.

It is often overlooked that one of the major advantages of NIR spectroscopy, i.e. its analysis speed and simplicity, provides an excellent opportunity to measure significantly more dosage units than what would be possible with the reference technique, i.e. HPLC. This is especially the case when automated NIR equipment is available. Hence, a massive amount of spectral information can be obtained easily to allow a better understanding of the possible spectral perturbations. The aim of the present study is to investigate whether the use of such spectral information on intra- and inter-batch variation (without corresponding reference analyses) during model development can improve the NIR model performance for predicting the API content in tablets from novel batches. The prediction performance of different approaches using the prior spectral information was compared to that of PLS models. It was also investigated whether these strategies allowed reducing the variability in the calibration set without compromising the prediction performance. This should lead to the definition of a cost-effective strategy for developing robust calibration models for the routine prediction of API in tablets.

## **2. Theory**

There are different ways to incorporate prior information into the calibration model. One distinguishes augmentation and orthogonal correction methods (supervised and unsupervised), while combinations of these are possible as well.

### **2.1. Augmentation methods**

#### *2.1.1. Noise augmentation in PLS*

The idea of noise augmentation (NA) or ensemble methods [10-13] is to expand the calibration set artificially in order to build PLS models that are more robust to the different kinds of

variability expected in the future sample population. An artificial calibration set can be created by augmenting the original one with a high number of spectral signatures, representing the 'noise'. The original spectra in the resulting calibration set should span the chemical variation in the best possible way, whereas the perturbation spectra should represent all the other possible variations that can be expected. One way to estimate the latter is through the measurement of spectra under varying perturbation conditions. Because the 'noise' added to the calibration spectra should be independent from the component of interest, this operation should not change the corresponding reference values.

### 2.1.2. Prediction – Augmented Classical Least Squares (P-ACLS)

In contrast to the inverse modeling approaches (e.g. PLS), classical least squares (CLS) regression is based on an explicit linear additive model (e.g. Lambert-Beer's law in spectroscopy). Equation (1) depicts the CLS model.

$$\mathbf{J} = \mathbf{PK} + \mathbf{E}_A \quad (1)$$

where  $\mathbf{J}$  is the matrix of the measured intensities,  $\mathbf{P}$  is the matrix of concentration values,  $\mathbf{K}$  is the matrix of the pure component signals at unit concentration, and  $\mathbf{E}_A$  the model error.

The major weakness of CLS is that it requires quantitative knowledge of all the spectrally active components in the calibration set to get an estimate of  $\mathbf{K}$ . Augmented (A)-CLS attempts to obtain a better estimation of  $\mathbf{K}$  by augmenting the predicted pure component matrix with empirically derived spectral shapes, e.g. the loading vectors from PCA on the spectral residuals  $\mathbf{E}_A$  (SRACLS), known pure component spectra [14], or a priori known other variation not included in the calibration set (PACLS) [15-16]. The addition of the spectral shapes both changes and corrects the concentration estimates for the component of interest.

## 2.2. Orthogonal correction methods

### 2.2.1. Orthogonal projections

There are essentially 2 contributions within the calibration matrix  $\mathbf{X}$ , i.e. one originating from the analyte of interest  $k$  ( $\mathbf{X}_k$ ), and another from all other sources of variance ( $\mathbf{X}_{-k}$ ).

$$\mathbf{X} = \mathbf{X}_k + \mathbf{X}_{-k} \quad (2)$$

Pretreatments based on orthogonal projections aim at removing those spectral patterns, which are 'interfering' with the desired prediction from the data matrix  $\mathbf{X}$ , before calibration on  $\mathbf{y}$ . It is attempted to return the spectra of  $\mathbf{X}$  as  $\mathbf{X}^*$ , containing the most condensed spectral information, i.e. the net analyte signal ( $\mathbf{X}_k$ ) [17]. The information to be removed, i.e. the clutter, can be defined based on pure component spectra of known interferences [18], an experimental design with varying perturbation factor(s) [19], or an (augmented) calibration data set [12]. The column vectors of the so obtained matrix  $\mathbf{S}$  form a basis of the detrimental subspace. Then, an orthogonal projector to  $\mathbf{S}$ , i.e.  $P_S^\perp$ , can be calculated to correct the initial matrix  $\mathbf{X}$  as follows.

$$\mathbf{X}^* = \mathbf{X}P_S^\perp = \mathbf{X}(\mathbf{I} - \mathbf{S}(\mathbf{S}^T\mathbf{S})^{-1}\mathbf{S}^T) \quad (3)$$

#### 2.2.1.1. External Parameter Orthogonalisation (EPO)

This unsupervised orthogonalisation method uses pure spectra (without reference value) to define a basis of the space spanned by the 'interfering' factors, this way estimating the parasitic subspace  $\mathbf{X}_k$  [19]. PCA is applied to  $\mathbf{D}$ , containing spectra collected while the perturbation factor is varying. Retaining only the first  $g_a$  PC's, the column vectors of the matrix of eigenvectors  $\mathbf{G}$  will represent an orthonormal basis of the subspace to be removed. Subsequently, an orthogonal projection is defined to filter the calibration spectra  $\mathbf{X}$  to obtain the 'corrected' ones ( $\mathbf{X}^*$ ).

$$\mathbf{X}^* = \mathbf{X}(\mathbf{I} - \mathbf{G}\mathbf{G}^T) \quad (4)$$

where  $\mathbf{G}$  is a matrix comprising the  $g_a$  first eigenvectors of the square matrix  $[\mathbf{D}^T\mathbf{D}]$ .

Because the original calibration data base  $\mathbf{X}$  is adjusted by means of orthogonal projection, the correction is embedded into the model. Hence, there is no need to reapply the correction to new spectra when using the model. Advantages of this method are the high flexibility and the fact that it does not require corresponding  $y$  values. Because it is based on external a priori information, the success of the unsupervised orthogonalisation will strongly depend on how one is able to identify the detrimental spaces within the variables space, without interfering with the useful space. It will depend on the comprehensiveness of the a priori spectral information and on the empirical tuning for identifying the information to be removed. In practice, the latter can be supported by assessing the RMSECV as a function of  $g_a$  and the number of latent variables, or by the evolution of the Wilks lambda value [19]. Nevertheless, using unsupervised orthogonal projection methods holds a risk of

removing too much information, i.e. when there is detrimental information not independent with the net analyte signal [20].

#### 2.2.1.2. Orthogonal Signal Correction (OSC)

OSC intends to subtract those factors from the calibration matrix  $\mathbf{X}$  which capture the variability in  $\mathbf{X}$  orthogonal to  $\mathbf{y}$  (supervised orthogonalisation) [21]. It uses a PLS-NIPALS-like algorithm, where the weighted regression vector  $\mathbf{w}$  is adjusted. Thus, OSC pre-processing involves two major steps, i.e. (1) estimation of the vector pair  $(\mathbf{t}, \mathbf{w})$  for which  $\mathbf{t}\mathbf{w}^T$  explains maximum variance in  $\mathbf{X}$  orthogonal to  $\mathbf{y}$ , and (2) removal of the contribution of the identified component from  $\mathbf{X}$ .

$$\mathbf{X}^* = (\mathbf{X} - \mathbf{t}\mathbf{p}^T) \quad (5)$$

Then, the newly obtained  $\mathbf{X}^*$  can be calibrated on  $\mathbf{y}$  by means of PLS. Different approaches have been proposed for estimating the score vector  $\mathbf{t}$ . As the orthogonal projection is not embedded in the calibration model, it has to be applied to new spectra prior to applying the calibration model [17].

#### 2.2.1.3. Net Analyte Pre-processing (NAP)

This method was introduced by Goicoechea and Olivieri [22]. To obtain an estimate of the parasitic subspace  $\mathbf{X}_{-k}$ , the following orthogonal projection of  $\mathbf{X}$  to  $\mathbf{y}$  is performed.

$$\mathbf{X}_{-k est} = (\mathbf{I} - \mathbf{y}(\mathbf{y}^T \mathbf{y})^{-1} \mathbf{y}^T) \mathbf{X} \quad (6)$$

In a next step,  $\mathbf{X}$  is projected orthogonal to the matrix  $\mathbf{U}$ , consisting of the first  $g_a$  PC's of  $\mathbf{X}_{-k est}$  to obtain  $\mathbf{X}^*$  for calibration on  $\mathbf{y}$  by means of PLS or CLS.

$$\mathbf{X}^* = \mathbf{X}(\mathbf{I} - \mathbf{U}\mathbf{U}^T) \quad (7)$$

The same transformation is applied to new spectra  $x$  prior to applying the calibration model.

$$x^* = (\mathbf{I} - \mathbf{U}\mathbf{U}^T)x \quad (8)$$

Thus, NAP has in common with OSC that it attempts removing the spectral information in  $\mathbf{X}$  that is orthogonal to  $\mathbf{y}$  (supervised orthogonalisation methods). Yet both methods use different routes to reach this goal. Compared to unsupervised orthogonalisation methods, OSC and NAP allow removing

the detrimental subspace in a less empirical way, and may also work efficiently when detrimental and useful information are not independent [20]. A disadvantage is that all the variation to be removed should be present in the calibration base, thus corresponding y values are necessary. Because of their close relationship to PLS, it has been repeatedly demonstrated that these methods can reduce the complexity of the PLS regression model, but do not enhance its predictive power [22-24].

### 3. Experimental

#### 3.1. NIR spectroscopy

Tablets were analyzed using a FT-NIR spectrometer (MPA, Bruker, Ettlingen, Germany). Spectra were collected in the 10000-5700  $\text{cm}^{-1}$  region with a resolution of 8  $\text{cm}^{-1}$  and averaged over 16 scans. The effective sample size was approximately 11% of the tablet. It was assumed that the API was uniformly distributed in the tablet, which was confirmed by studies during method development.

#### 3.2. HPLC

HPLC was used as the reference technique. All HPLC measurements were performed at Janssen Pharmaceutica, Beerse, Belgium, using their developed and validated method (confidential).

#### 3.3. Data analysis

PCA, PLS regression and OSC pre-treatment were performed in Matlab 7.1. (The Mathworks, Nattick, MA) using the PLS toolbox 6.2. (Eigenvector Research, Wenatchee, WA). EPO, ACLS and NAP were directly programmed in Matlab. Minitab 16 (Minitab, PA) was used for testing the statistical significance of the predictions by the different models.

##### 3.3.1. Model development and validation

The calibration set consisted of oblong shaped tablets of 6 target API concentration levels, i.e. 0, 25, 50, 100, 150 and 187.5 mg API with respective 0, 7.28, 14.56, 29.12, 38.83 and 48.54%  $\text{ww}^{-1}$  per tablet. Tablets at the extreme concentration levels originated from laboratory batches, and were manufactured to extend the concentration range spanning the chemical variability within the specification limits. The tablets at the other concentration levels were manufactured in a production line and represent tablets to be marketed. As the calibration set contained a large concentration range, the tablet composition also varied over the concentration levels.



The NIR spectra were limited to the 9000-7600  $\text{cm}^{-1}$  range and were pre-processed with standard normal variates (SNV). Except for the ACLS and NAP-CLS models, all other calibration models were built with PLS. Mean-centering was always performed prior to PLS modelling. The optimal number of latent variables for each model was selected based on the minimal RMSECV from 'leave-one-batch-out' cross validation. Two independent test sets were used for external model validation. Test set 1 consisted of 54 tablets, originating from 4 new production batches (one for each dose to be marketed) and 2 laboratory batches at the extreme concentration levels. In test set 2 (30 tablets from 4 production batches) strong water variations were introduced by storing the tablets at different conditions (see further in 3.3.2.). The root mean square error of prediction (RMSEP) was used as the performance criterion to assess and compare the predictive abilities of the different models. The significance of differences in prediction power was assessed by a two-way ANOVA, performed on the absolute values of the prediction errors [25]. Multiple comparisons were performed by the Dunnett's test, using the PLS model as a control to calculate  $p$  values.

### 3.3.2. Available prior information

NIR spectra were recorded from the pure tablet components, i.e. from the API (**A**) and the excipients (**B**). Fig.1 shows the obtained spectra in the considered spectral range after SNV pre-processing. 10086 NIR spectra from tablets originating from 27 different production batches (i.e. for doses to be marketed: 3 batches at 7.28 % w/w, 12 batches at 14.56 %w/w, 7 batches at 29.12 % w/w and 5 batches at 38.83 % w/w) were measured. They included tablets produced during pharmaceutical development and clinical trials. They also included characterization batches, produced at the limits of the critical process settings, to increase the spectral variability originating from process variability. Since hundreds of tablets per batch were measured, systematic variations in the spectra due to process variation over time may also be covered. The NIR spectra of the tablets originating from laboratory batches (at the extreme target API concentration levels) were added to this matrix, creating **C** (10213 x 364). 107 spectra from processed tablets with deliberate strong water variations were measured (Fig.2). Either the tablets were stored for 30 min at 50°C (to decrease the water content in the tablets), either for 3h at 75% RH, or for 16h at 75% RH (to increase the water content in the tablets).

## 3.4. Strategies for efficient use of the prior information

### 3.4.1. Identifying the information to be removed (clutter)

The strategies applied in this paper do not require exact knowledge of all perturbation sources, but just need a representative estimate of the spectral variation to be expected in tablets

from future batches. This can be obtained by making use of the extremely fast and simple NIR analysis, where far more tablets could be analyzed than what would be possible with HPLC. Fig.3 shows a PCA score plot (PC1 versus PC2) of samples from randomly selected batches (coloured markers) for calibration of the initial PLS model, and all the prior spectral information **C** (10213 x 364) are in grey dots. The scores for **C** indicate a generally higher spread for the production batches, and this was confirmed for the following investigated PC's as well (not shown). The PCA score plots of the spectra of **C** from tablets at each target concentration level also showed clear systematic effects between and within production batches. Examples are shown in Figs. 4-5. An interesting example of intra-batch variability is shown in Fig.4, where the spectra from the batch marked with blue squares showed two separate clusters. This may be due to a systematic effect that occurred over the manufacturing time of the batch. Hence, selecting the samples for calibration holds a risk of not fully covering the spectral intra- and inter-batch variation that may be present in future samples.

As the variation in **C** is partly due to differences in API content (the parameter to be quantified) in each individual tablet, this spectral information cannot be used 'as such' for defining the clutter for the proposed strategies. Many correction methods define the clutter by the difference spectra at various perturbation levels, which works well when the property to be determined can be held constant [17]. Here, it is impossible to obtain individual tablets, even from the same batch, with exactly the same API concentration and tablet composition. Another possibility is using a y-gradient method to select the spectra for calculating the difference spectra [26], but no y-reference values are available for the massive amount of prior spectral information. For the present case study, we proposed another methodology to obtain a relevant estimate of the clutter. To correct **C** for known chemical variations, the pure spectra from the tablet constituents were used to calculate a basis. Removal of their spectral contributions from **C** was performed by orthogonal projection [18]. After the correction of **C**, the spectra have lost their initial form, leaving only the spectral shapes from other variations (Fig.6).

Fig.7 schematically illustrates the different tested strategies for obtaining a clutter. Matrix **B** (tablet excipients) can be defined as a clutter containing the known chemical noise. Considering both **A** and **B** for the orthogonal correction renders a clutter containing only the physical variations and unknown chemical variation (e.g. water content...) (**D**). Using **A** yields a clutter containing physical variations and chemical variability caused by all 'interfering' excipients (**E**).

### 3.4.2. Strategies using EPO

Depending on how the clutter is defined, there are different possibilities to perform EPO corrections on the calibration spectra **X** (Fig.7). In a similar way as above,  $EPO_{chem}$  directly removes the known interfering chemical variations (**B**) from the **X** spectra. In  $EPO_{glob}$  the expected chemical

1 and other variations were captured in  $\mathbf{E}$  for global EPO correction of  $\mathbf{X}$ . Here, optimization of the  
2 dimension of the subspace to be withdrawn was performed by analyzing the evolution of the  
3 RMSECV as a function of  $g_o$  and the number of LV's [19].  
4

### 5 3.4.3. Strategies using data augmentation

6 Noise augmentation (NA) of the calibration set was obtained by adding the mean-centered  
7 spectra of  $\mathbf{D}$  to  $n$  repetitions of the original calibration spectra  $\mathbf{X}$ . Similarly,  $\mathbf{y}_{aug}$  was developed as  $n$   
8 repetitions of  $\mathbf{y}$ , because the added spectral variations are not supposed to change the  
9 corresponding  $y$  value of each spectrum (Fig.7). Hence, a calibration set containing 10213 objects was  
10 obtained and used for PLS regression.

11 Because OSC and NAP need a calibration data base to calculate the orthogonal projection,  
12 these techniques were applied on the augmented calibration set, i.e.,  $\mathbf{X}_{aug}$  and  $\mathbf{y}_{aug}$  containing the  
13 most complete information. The  $OSC_{aug}$  used the  $\mathbf{X}_{aug}$  and  $\mathbf{y}_{aug}$  calibration base for calculating the  
14 orthogonal correction. The NA-NAP strategies performed NAP on the augmented calibration base  
15  $\mathbf{X}_{aug}$  and  $\mathbf{y}_{aug}$  this way calculating a projection factor to transform the spectra in  $\mathbf{X}$  to  $\mathbf{X}^*$  prior to PLS or  
16 CLS regression. The optimal number of factors  $g_o$  to be removed was determined via the minimal  
17 RMSECV. The test set spectra underwent the same transformation as the calibration spectra.  
18

### 19 3.4.4. Strategies using ACLS

20 The SRACLS procedure described in [14] was used to define an augmented pure component  
21 matrix  $\mathbf{K}$  consisting of the pure spectra from the API ( $\mathbf{A}$ ) and the excipients ( $\mathbf{B}$ ), and the  $p$  first loading  
22 vectors from the residual matrix  $\mathbf{E}_A$ . In the PACLS procedure the augmented pure component matrix  
23 of the SRACLS procedure was further augmented with the  $k$  first loading vectors from  $\mathbf{D}$ . The optimal  
24  $p$  and  $k$  were determined through cross-validation (minimal RMSECV).  
25

## 26 3.5. Model update for strong water variations

27 To assess whether the strategies using prior spectral information can account for strong  
28 water variations, 2 different model update strategies were proposed. In *model update 1*, the spectra  
29 of such processed tablets were added (with corresponding reference values) to the calibration set 1  
30 to obtain an updated calibration set containing 197 samples (120 initial + 77 processed samples). The  
31 prior information contained, just as in the previous section, no spectra of processed tablets. In *model*  
32 *update 2*, the initial training set 1 (120 not-processed tablets) was used for calibration. The prior  
33 information matrix containing intra- and inter-batch variability was augmented with spectra from  
34 processed tablets to account also for the strong water variations.  
35

## 4. Results and discussion

### 4.1. Models using prior information containing intra- and inter-batch variation

In this section, the above discussed strategies were applied to evaluate their capability for filtering the intra- and inter-batch effects hampering the model performance. Test set 1 was used to evaluate the model's ability to withstand unknown variability (external validation). With PLS models covering different amounts of variation in their calibration set as a benchmark, the performance of the investigated modelling strategies was investigated in terms of number of latent variables (LV's), root mean square error of prediction (RMSEP) and its statistical significance (Table 1). The number of factors used for the correction method was also reported. All models were trained on 4 different sample sets, holding different amounts of variation. Set 1 (120 tablets, 12 batches) contained the most intra- and inter-batch variation, while set 4 (18 tablets, 6 batches) contained the least. Sets 2 (60 tablets, 6 batches) and 3 (30 tablets, 6 batches) contained more intra-batch variation compared to set 4. The prior information matrix **C** was considered to be representative for intra- and inter-batch variability in the spectra.

The PLS, NAP-CLS, NAP-PLS and OSC-PLS models did not use any prior information and were added for comparison purposes. They were marked in Table 1 with an asterisk. For all calibration sets the OSC-PLS, NAP-PLS and NAP-CLS models were generally more parsimonious (less LV's) compared to PLS, but did not show improved predictive power. Many papers present similar conclusions on this behavior of OSC and NAP [22-24]. PLS established higher prediction errors in models calibrated on sample sets containing less variability.

The prediction performances of the ACLS models (SRALCS and PACLS) were much worse than the PLS-based models, making the latter models more suitable for the present case study. It is hypothesized that ACLS models, based on a linear-additive model may have trouble capturing the physical differences between the NIR spectra of powders (pure components) and tablets. They might be more suitable for modelling pure mixtures of different components [14]. Augmenting the calibration set prior to PLS, i.e. by noise addition (NA)-PLS, yielded models with more LV's compared to the original PLS models. This can be attributed to more spectral variation being included in the calibration set. Noise augmentation before PLS regression (NA-PLS), making the calibration matrix very large (10213 objects) and noisy, could not significantly enhance the prediction performance either, nor could OSC when applied on this augmented calibration matrix (NA-OSC-PLS). NA-NAP-PLS showed a statistically better prediction performance than the PLS model ( $p < 0.05$ ), with a 34-40% lower RMSEP, while needing less latent variables. A similar behavior of this approach was also seen on all other calibration sets with a progressively smaller size (Table 1).

1 Compared to NA-PLS, adding the NAP pre-processing seems to be advantageous for making  
2 the model more parsimonious and for improving the prediction performance. Nadler and Coifman  
3 [27] concluded that PLS may behave differently in two idealized settings: for a noise-free training set  
4 the regression vector computed by PLS is, up to normalization, the net analyte signal, while for a  
5 noisy infinite training set the regression vector is not purely proportional to the NAS vector, but is  
6 optimal under a mean squared error of prediction criterion. Considering the obtained results, it may  
7 be speculated that NAP may have a different relationship towards PLS in the augmented calibration  
8 set. It would be an interesting topic of further study to find out why the NAP algorithm performs so  
9 much better on the large and noisy augmented calibration sets, and actually allows here an  
10 improvement in prediction performance.

11 Compared to NAP-PLS, the NA procedure may add more variability into the calibration data  
12 base, which is missing when using non-exhaustive calibration. In the NAP-CLS and NAP-PLS models  
13 only 7 to 8 factors were removed by the NAP procedure, whereas the NA-NAP operation removed  
14 between 10 and 12 factors, and approximately 1-3 additional factors were needed for subsequent  
15 PLS regression. This may approximate the number of perturbation factors to be expected in future  
16 tablets. Hence the augmentation step is also necessary to make the model more robust and the  
17 completeness of the clutter is an important factor for the performance and robustness of the  
18 models.

19 Except for the model calibrated on set 4 the NA-NAP method performed clearly better when  
20 combined with PLS instead of CLS. In these cases, PLS was able to further model the important  
21 information. In case of calibration set 4, the NA-NAP operation returned the original calibration  
22 spectra  $\mathbf{X}$  into  $\mathbf{X}^*$  containing almost exclusively information on the component to be quantified  
23 (Fig.8), with PLS and CLS giving a comparable result.

24 Compared to PLS, EPO-PLS based calibrations showed lower RMSEP's and needed a lower  
25 number of LV's, and especially for models calibrated on sample sets carrying fewer variability there  
26 was a statistical difference between both approaches. Considering only the chemical variations in the  
27 correction, i.e.  $EPO_{chem}$ , allowed improving the model performance to some extent. This can be  
28 attributed to the correction for the changing tablet compositions over the API target concentration  
29 levels, as well as for excipient inhomogeneities among individual tablets.

#### 30 31 4.2. Deliberate water variation (model update)

32 When the tablets are processed, i.e. stored at different extreme temperature and humidity  
33 conditions, the water content in the tablets can change considerably [6]. Compared to the tablets  
34 that were stored under normal conditions, those stored long in high humidity conditions (i.e. 75%  
35 RH) showed an interesting effect in their NIR spectra, i.e. shifts in the API absorption signals were

visible (Fig.2). These frequency shifts may suggest a different hydrogen bonding state of the API [28]. As this variation affects the PLS model predictions (Table 2), we investigated whether this information might be considered in the calibration base or in the clutter for a model update.

In case the model was updated by adding deliberate water variation to the calibration data base, the NA-NAP-PLS strategy again outperformed the others, and was the only model being statistically significantly better than PLS (Table 2, update 1). With a model requiring only 4 LV's an RMSEP of 0,380 for test set 1 and one of 0,385 for test set 2 were obtained. Hence the model was able to predict perturbed samples reasonably well (test set 2), while also not substantially reducing the prediction performance for normal or unperturbed samples (test set 1). Compared to NA-NAP-CLS, NA-NAP-PLS allowed a lower RMSEP for both test sets. The NA-NAP-PLS model now required 1 LV more than before model update (Table 1, set 1), which it probably needs to model the shifting of the API absorption signals.

The second model updating strategy included the spectra of the perturbed samples in the prior spectral information matrix. As expected, models using no prior information from processed tablets, i.e. PLS, NAP-CLS, NAP-PLS, OSC-PLS, SRACLS and  $EPO_{chem}$ -PLS, displayed poor prediction performance for test set 2 (Table 2, update 2). This could be attributed to the fact that they did not take into account the deliberate water variation. Although some of the modelling strategies that used prior information from processed tablets were able to get better prediction errors for test set 2, the overall model performance was low and none of the methods performed significantly better than PLS. Compared to the models updated with strategy 1, none of the methods in model update 2 was able to obtain good prediction results for test set 2. The prediction performance for test set 1 remained reasonably well in all models (except for PACLS), whereas predictions for test set 2 were always poor. The reason for this was found to be that the spectral perturbations due to strong water content variations in C (also causing shifting of the API absorption signals) are not completely independent of  $\mathbf{y}$ , and therefore could not be considered in the model update. The proposed method to define the clutter (step 1 of the strategy, see section 4.1.1) uses unsupervised orthogonal projections with the pure component spectra. These may remove that part of  $\mathbf{C}$ , which is co-linear to  $\mathbf{y}$ , making it not possible to fully consider this kind of spectral variation in the clutter. These results demonstrate that the proposed strategies only allow making the model more robust for spectral variation that is not co-linear with the one from the property of interest. Hence, this variability should still be considered in the calibration set.

## 5. Conclusions

When developing a NIR model for routine quality control that needs to operate over a longer time course, it may not be evident to cover all possible perturbation factors in the calibration set.

1 Because it is difficult to select samples describing intra-and inter-batch variability and to decide when  
2 the calibration set is representative, this study explored whether the judicious use of such prior  
3 spectral information during model development can improve the performance of NIR models for  
4 predicting the API content in tablets. The proposed strategies did neither require exact knowledge of  
5 the perturbation levels (only controlled variability of the considered perturbations was necessary),  
6 nor needed extra reference analyses.

7 Rather than exhaustive calibration, a more cost-effective model development approach was  
8 aspired. A massive amount of prior spectral information on intra- and inter-batch variation was  
9 obtained to allow a more representative view of the possible disturbing (i.e. systematic) effects to be  
10 encountered in the NIR spectra of tablets from a routine production. The best approach (NA-NAP-  
11 PLS) consists of four essential steps. First, the disturbances to be removed are identified by means of  
12 using an orthogonal projection of pure component spectra (API and excipients) to those of the  
13 tablets containing representative intra- and inter-batch variability. This way, the known chemical  
14 variation is removed from the spectra leaving only the other variations. These mean-centered  
15 variations are in a second step added to  $n$  repetitions of the spectra of  $\mathbf{X}$  for calibration (noise  
16 augmentation) to increase the variability in the existing calibration database. In a third step, net  
17 analyte pre-processing is used to remove those variations in  $\mathbf{X}_{aug}$  that are orthogonal to  $\mathbf{y}_{aug}$ , in order  
18 to obtain a better estimate of the NAS. Finally, in a fourth step the orthogonally corrected  $\mathbf{X}$  spectra  
19 are regressed to  $\mathbf{y}$  by PLS regression.

20 Compared to PLS, a statistically significant improvement in prediction performance and a 34-  
21 40% reduction in RMSEP was obtained for predicting new tablets with unknown intra- and inter-  
22 batch variability. The model needed a minimal amount of calibration samples and latent variables. It  
23 would be interesting to further study the behavior of NAP-PLS in very large and noisy calibration  
24 matrices.

25 Where PLS models need the calibration set to be as complete as possible, models using prior  
26 spectral information require completeness of their clutter. This was also demonstrated by updating  
27 the model for strong water content variations in the tablets. The latter caused frequency shifts of the  
28 API absorption signals that could not be included in the clutter, resulting in reduced performance of  
29 the proposed strategies for this kind of variation. It also revealed a limitation of the method: the  
30 spectral variance that is co-linear with the wanted variation could not be considered in the clutter  
31 and should be considered in the calibration set.

## Acknowledgements

Sigrid Pieters and Wouter Saeys are respectively funded as aspirant and postdoctoral fellow of the Research Foundation – Flanders (FWO).

## References

- [1] J. Luypaert, D.L. Massart, Y. Vander Heyden, *Talanta* 72 (2007) 865-883.
- [2] T. De Beer, A. Burggraeve, M. Fonteyne, L. Saerens, J.P. Remon, C. Vervaet, *Int. J. Pharm.* 147 (2011) 32-47.
- [3] D. Cozzolino, W. Cynkar, N. Shah, P. Smith, *Anal. Bioanal. Chem.* 401 (2011) 1475-1484.
- [4] W.F. McClure, *Anal. Chem.* 66 (1994), 43-53.
- [5] M. Blanco, A. Peguero, *J. Pharm. Biom. Anal.* 52 (2010) 59-65.
- [6] T. Van den Kerkhof, R. De Maesschalck, K. Vanhoutte, M.C. Coene, *J. Pharm. Biomed. Anal.* 42 (2006) 517-522.
- [7] A. Peirs, J. Tirry, B. Verlinden, P. Darius, B.M. Nicolaï, *Postharv. Biol. Technol.* 28 (2003) 269-280.
- [8] EMEA/CHMP/CVMP/QWP/17760/2009 Rev 2 (2012), accessed on 28/03/2012
- [9] D.L. Massart, B.M.G. Vandeginste, L.M.C. Buydens, S. De Jong, P.J. Lewi, J. Smeyers-Verbeke, *Handbook of chemometrics and qualimetrics: part A*, Elsevier, Amsterdam, 1997.
- [10] M.J. Saiz-Abajo, B.H. Mevik, V.H. Segtnan, T. Naes, *Anal. Chim. Acta* 533 (2005) 147-159.
- [11] B.H. Mevik, V.H. Segtnan, T. Naes, *J. Chemom.* 18 (2004) 498-507.
- [12] J.A. Fernandez Pierna, F. Chauchard, S. Preys, J.M. Roger, O. Galtier, V. Baeten, P. Dardenne, *Chemom. Intell. Lab. Syst.* 106 (2011) 152-159.
- [13] V.H. Segtnan, B.H. Mevik, T. Isaksson, T. Naes, *Appl. Spectr.* 59 (2005) 816-825.
- [14] W. Saeys, K. Beullens, J. Lammertyn, H. Ramon, T. Naes, *Anal. Chem.* 80 (2008) 4951-4959.
- [15] D. M. Haaland, D.K. Melgaard, *Appl. Spectr.* 54 (2000) 1303-1312.
- [16] D.M. Haaland, D.K. Melgaard, *Vib.Spec.* 29 (2002) 171-175.
- [17] J.C. Boulet, J.M. Roger, *Chemom. Intell. Lab. Syst.*, in press
- [18] J.C. Boulet, T. Doco, J.M. Roger, *Chemom. Intell. Lab. Syst.* 87 (2007) 295-302.
- [19] J.M. Roger, F. Chauchard, V. Bellon-Maurel, *Chemom. Intell. Lab. Syst.* 66 (2003) 191-204.
- [20] S. Preys, J.M. Roger, J.C. Boulet, *Chemom. Intell. Lab. Syst.* 91 (2008) 28-33.



- [21] S. Wold, H. Antti, F. Lindgren, J. Ohman, *Chemom. Intell. Lab. Syst.* 44 (1998) 175-185.
- [22] H. C. Goicoechea, A. C. Olivieri, *Chemom. Intell. Lab. Syst.* 56 (2001) 73-81.
- [23] W. Ni, S.D. Brown, R. Man, *Chemom. Intell. Lab. Syst.* 98 (2009) 97-107.
- [24] O. Svensson, T. Kourti, J.F. MacGregor, *J. Chemom.* 16 (2002) 176-188.
- [25] H.R. Cederkvist, A.H. Aastveit, T. Naes, *J. Chemom.* 19 (2005) 500-509.
- [26] [http://wiki.eigenvector.com/index.php?title=Declutter Settings Window](http://wiki.eigenvector.com/index.php?title=Declutter%20Settings%20Window), accessed on 11/04/2012.
- [27] B. Nadler, R.R. Coifman, *J. Chemom.* 19 (2005) 45-54.
- [28] S. Pieters, T. De Beer, J.C. Kasper, D. Boulpaep, O. Waszkiewicz, M. Goodarzi, C. Tistaert, W. Friess, J.P. Remon, C. Vervaet, Y. Vander Heyden, *Anal. Chem.* 84 (2012) 947-955.

## Figure captions

**Fig. 1.** NIR spectra of the pure tablet components (SNV pre-processed) in the range 9000-7600  $\text{cm}^{-1}$ .

**Fig. 2.** NIR spectra of processed and unprocessed tablets (SNV pre-processed) in the range 9000-7600  $\text{cm}^{-1}$ . Spectra from unprocessed tablets are in grey, from tablets stored for 30 min at 50°C in green, from tablets stored for 3h and 16h at 75% RH in purple and red, respectively. The black ellipses indicate visible shifting of the API absorption signals.

**Fig. 3.** PCA score plot (PC1 versus PC2) from **C** covering the 10213 spectra of tablets from 27 different batches (in grey). The spectra used in the calibration set 1 are represented by differently coloured markers according to their batch.

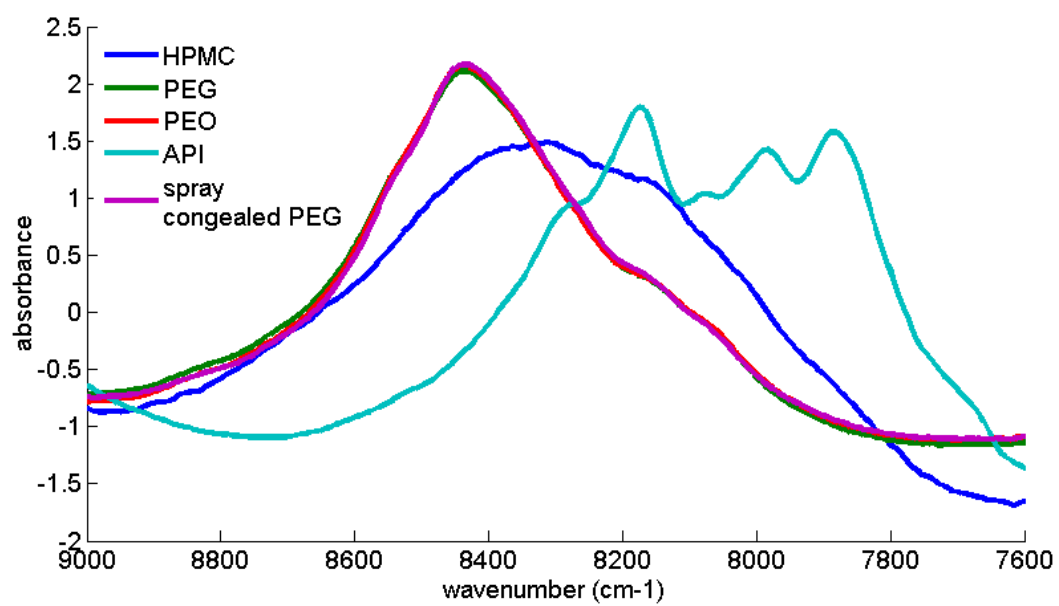
**Fig. 4.** PCA score plot (PC1-PC2-PC3) of NIR spectra (SNV pre-processed) from tablets with 25 mg target API concentration manufactured in a production line. The different batches are indicated with differently colored markers.

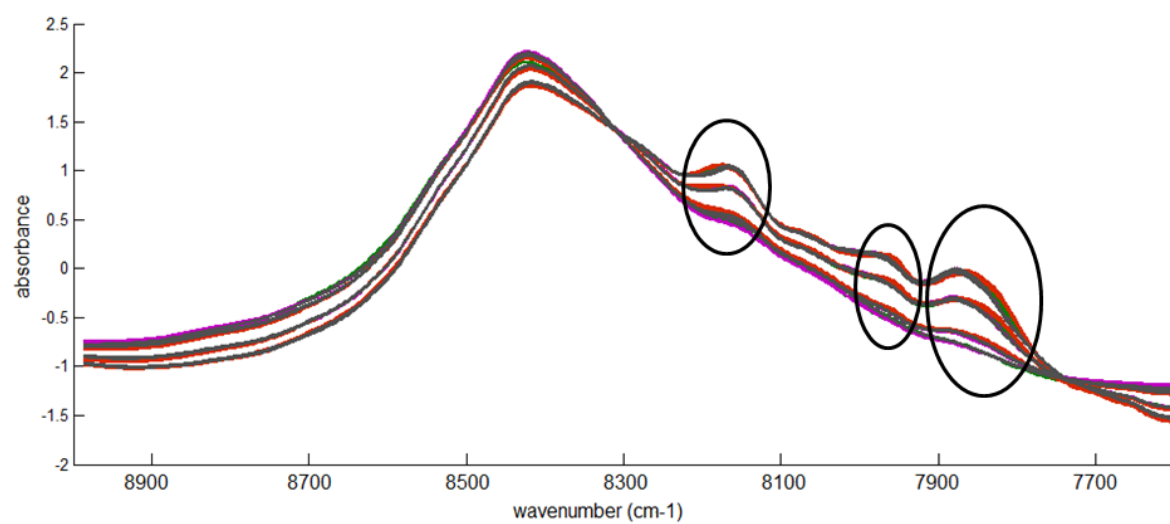
**Fig. 5.** PCA score plot (PC1-PC2-PC3) of NIR spectra (SNV pre-processed) from tablets with 100 mg target API concentration manufactured in a production line. The different batches are indicated with differently colored markers.

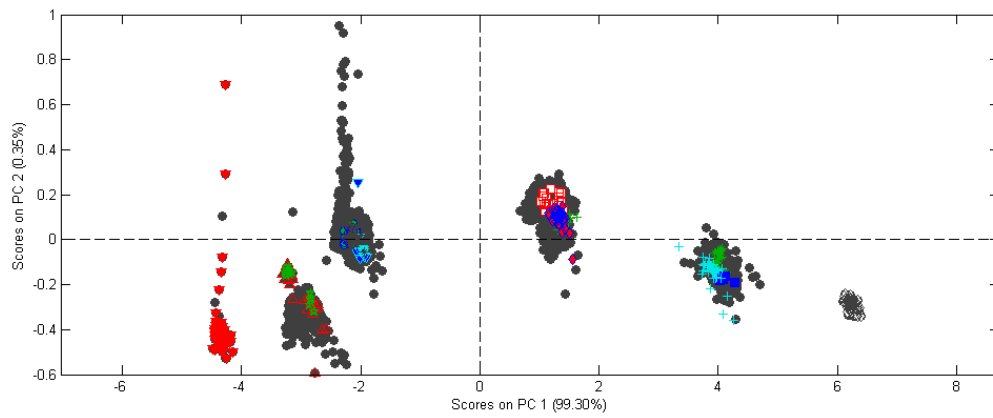
**Fig.6.** Spectra before (A) and after (B) orthogonal correction using the pure component spectra of the API and excipients to remove the known chemical variations.

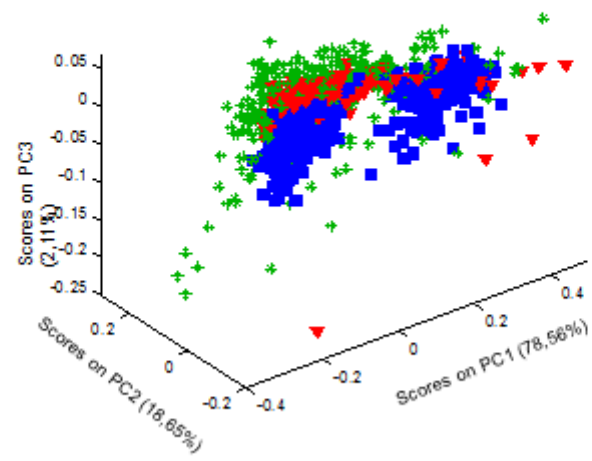
**Fig.7.** Flow chart depicting the applied methodologies. **A**, **B** and **C** are matrices containing prior information; the subscript *c* represents mean-centered matrices. **X** (measured spectral absorbances **J**) and **y** (concentrations to be predicted **P**) represent the original calibration set. **X**<sub>aug</sub> and **y**<sub>aug</sub> are the noise augmented calibration set.  $P^\perp$  represents orthogonal projection and **X**<sup>\*</sup> is the calibration matrix returned after correction. **K** is the matrix of pure component signals at unit concentration.

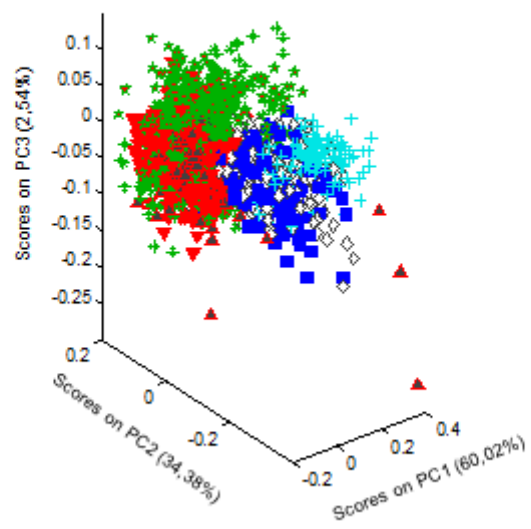
**Fig. 8.** Calibration spectra **X**<sub>nap</sub><sup>\*</sup> (18 x 364) of set 4 after NA-NAP transformation. These spectra were subsequently calibrated on **y** with PLS. The different colors represent the different API concentration levels.

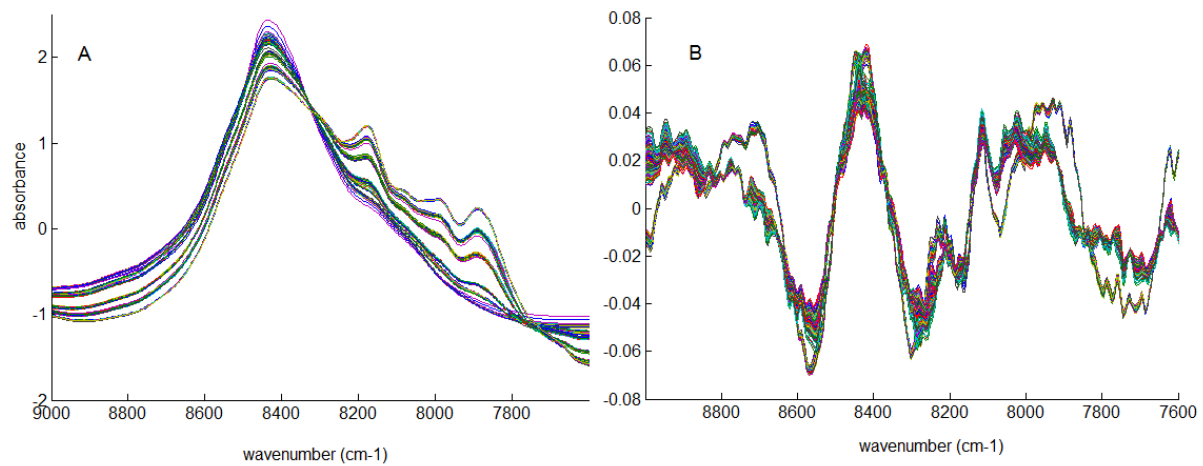




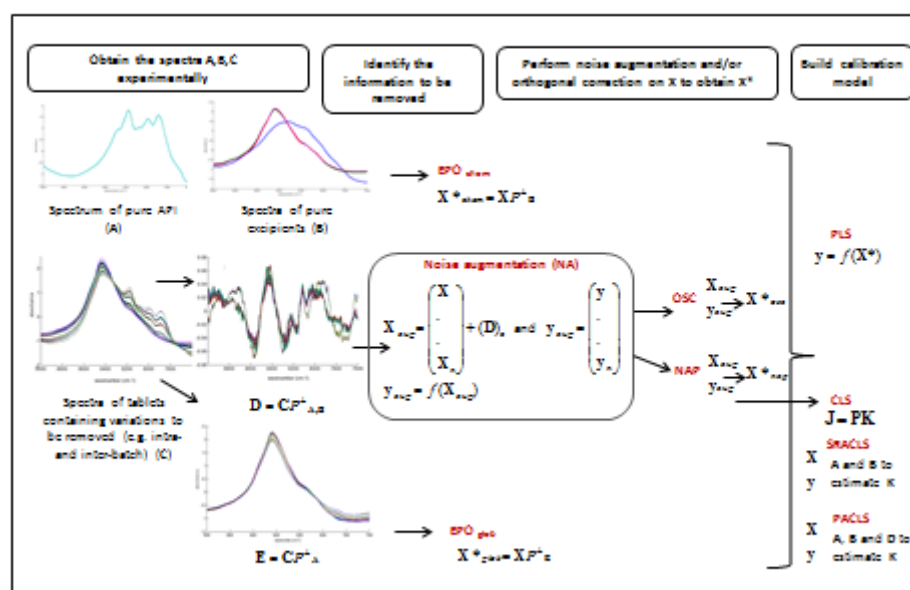


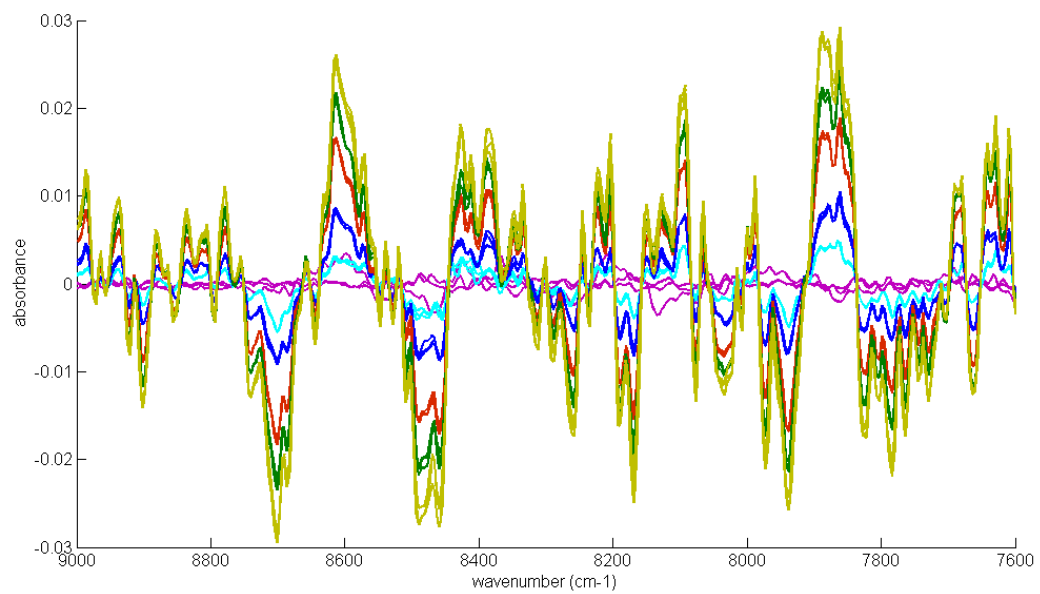












**Table 1.** Prediction performance (for test set 1) of the different models calibrated on sample sets with different variation. The matrix with prior intra- and inter-batch variation was **C** (10213 x 364). The best prediction results are in bold.

Calibration set	Set 1 (120 tablets, 12 batches)				Set 2 (60 tablets, 6 batches)				Set 3 (30 tablets, 6 batches)				Set 4 (18 tablets, 6 batches)			
Model	# factors	# L V	RMS EP	<i>p</i> value <sup>b</sup>	# factors	# L V	RMS EP	<i>p</i> value <sup>b</sup>	# factors	# L V	RMS EP	<i>p</i> value <sup>b</sup>	# factors	# L V	RMS EP	<i>p</i> value <sup>b</sup>
PLS <sup>a</sup>	-	7	0.58 6	-	-	4	0.63 0	-	-	4	0.63 6	-	-	4	0.66 3	-
NAP-CLS <sup>a</sup>	7	-	0.54 7	0.84 9	8	-	0.73 9	0.54 3	7	-	0.75 5	0.51 2	7	-	1.12 6	<0.0 01 <sup>c</sup>
NAP-PLS <sup>a</sup>	7	4	0.60 1	1.00 0	8	5	0.60 8	1.00 0	7	4	0.60 4	0.80 7	7	4	1.12 3	<0.0 01 <sup>c</sup>
OSC-PLS <sup>a</sup>	3	4	0.56 0	1.00 0	3	5	0.64 6	0.98 6	1	3	0.63 6	1.00 0	1	3	0.66 3	0.99 9
NA-PLS	-	8	0.49 6	0.92 8	-	5	0.58 8	1.00 0	-	6	0.56 6	0.96 3	-	6	0.58 5	0.98 5
SRACLS	5	-	0.65 9	0.99 4	10	-	1.25 4	<0.0 01 <sup>c</sup>	6	-	0.97 5	<0.0 01 <sup>c</sup>	6	-	0.98 8	<0.0 01 <sup>c</sup>
PACLS	5-1	-	0.90 8	<0.0 01 <sup>c</sup>	8-1	-	1.09 9	<0.0 01 <sup>c</sup>	6-1	-	0.94 5	<0.0 01 <sup>c</sup>	6-1	-	1.13 9	<0.0 01 <sup>c</sup>
EPO <sub>chem</sub> -PLS	4	5	0.51 1	0.71 6	4	4	0.46 8	0.01 3 <sup>c</sup>	4	5	0.48 7	0.10 7	4	5	0.51 9	0.28 1
EPO <sub>gr</sub> -PLS	1	6	0.46 1	0.44 9	1	7	0.58 1	1.00 0	1	6	0.47 8	0.00 4 <sup>c</sup>	1	7	0.51 6	0.00 9 <sup>c</sup>
NA-NAP-CLS	10	-	0.57 2	1.00 0	12	-	0.52 7	0.58 8	11	-	0.54 6	0.37 0	11	-	0.43 4	0.02 6 <sup>c</sup>
NA-NAP-PLS	<b>10</b>	<b>3</b>	<b>0.38 8</b>	<b>0.02 7<sup>c</sup></b>	<b>12</b>	<b>2</b>	<b>0.40 7</b>	<b>0.00 7<sup>c</sup></b>	<b>11</b>	<b>2</b>	<b>0.38 2</b>	<b>&lt;0.0 01<sup>c</sup></b>	<b>11</b>	<b>1</b>	<b>0.43 1</b>	<b>0.01 5<sup>c</sup></b>
NA-OSC-PLS	2	3	0.48 8	0.69 7	1	2	0.54 6	0.58 6	1	2	0.53 3	0.42 4	1	2	0.56 4	0.72 7

<sup>a</sup> These models did not use prior information; they used only information from the calibration set

<sup>b</sup> *p* values (Dunnett's test) for the significance testing by two-way ANOVA of the predictive ability compared to PLS

<sup>c</sup> Indicates significantly lower absolute prediction errors on significance level  $\alpha=0.05$

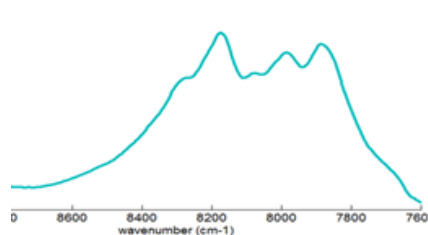
**Table 2.** Prediction performance (for test sets 1 and 2) after model update for deliberate water variations. The best prediction results are in bold.

Model	Update 1					Update 2				
	# factors	# LV	RMSEP test set 1	RMSEP test set 2	<i>p</i> value <sup>b</sup>	# factors	# LV	RMSEP test set 1	RMSEP test set 2	<i>p</i> value <sup>b</sup>
PLS <sup>a</sup>	-	8	0.632	0.407	-	-	7	<b>0.586</b>	1.274	-
NAP-CLS <sup>a</sup>	8	-	0.534	0.461	0.883	7	-	0.547	1.488	1.000
NAP-PLS <sup>a</sup>	8	6	0.644	0.302	0.645	7	4	0.601	1.358	1.000
OSC-PLS <sup>a</sup>	5	6	0.632	0.346	0.989	3	4	0.560	1.506	1.000
NA-PLS	-	8	0.548	0.433	0.549	-	5	0.466	0.854	0.844
SRACLS <sup>a</sup>	9	-	0.780	0.466	0.124	5	-	0.659	1.557	1.000
PACLS	9-1	-	0.756	0.417	0.067	5-12	-	5.490	8.323	<0.001 <sup>c</sup>
EPO <sub>chem</sub> -PLS <sup>a</sup>	4	5	0.575	0.376	0.792	4	5	0.511	0.913	0.948
EPO <sub>gl</sub> -PLS	1	7	0.516	0.397	0.284	1	7	0.549	0.977	0.989
NA-NAP-CLS	11	-	0.512	0.517	0.964	13	-	0.567	0.728	0.923
NA-NAP-PLS	<b>11</b>	<b>4</b>	<b>0.380</b>	<b>0.385</b>	<b>0.003<sup>c</sup></b>	13	3	0.488	0.761	0.595
NA-OSC-PLS	4	2	0.498	0.409	0.265	1	2	0.510	0.976	0.966

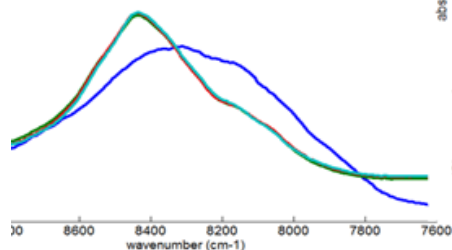
<sup>a</sup>These models did not use prior information of processed tablets containing deliberate water variations.<sup>b</sup>*p* values (Dunnett's test) for the significance testing by two-way ANOVA of the predictive ability compared to PLS.

<sup>c</sup>Indicates significantly lower absolute prediction errors on significance level  $\alpha=0.05$ .

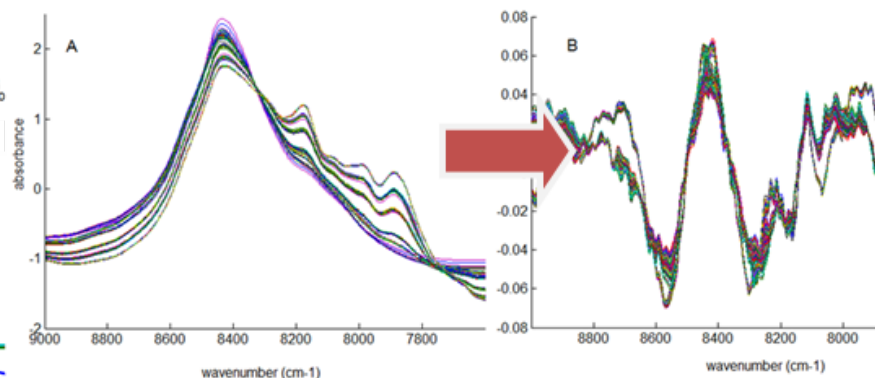
$$CP^{\perp}_{A, B}$$



(A) pure API



(B) pure excipients



(C) spectra from tablets with variations to be removed

(D) clutter

**Highlights:**

- a cost-effective NIR model development strategy is proposed for API content prediction in tablets
- judicious use of prior spectral information improves PLS model performance
- the clutter captures representative intra- and inter-batch spectral variability to be removed
- model requires completeness of the clutter rather than comprehensive calibration sets